# Optimizing Bloom Filters for Modern GPU Architectures

Daniel Jünger*
djuenger@nvidia.com
NVIDIA Corporation
Santa Clara, USA

Kevin Kristensen*
kckristensen@wisc.edu
University of Wisconsin-Madison
Madison, USA

Yunsong Wang
yunsongw@nvidia.com
NVIDIA Corporation
Santa Clara, USA

Xiangyao Yu
yxy@cs.wisc.edu
University of Wisconsin-Madison
Madison, USA

Bertil Schmidt
bertil.schmidt@uni-mainz.de
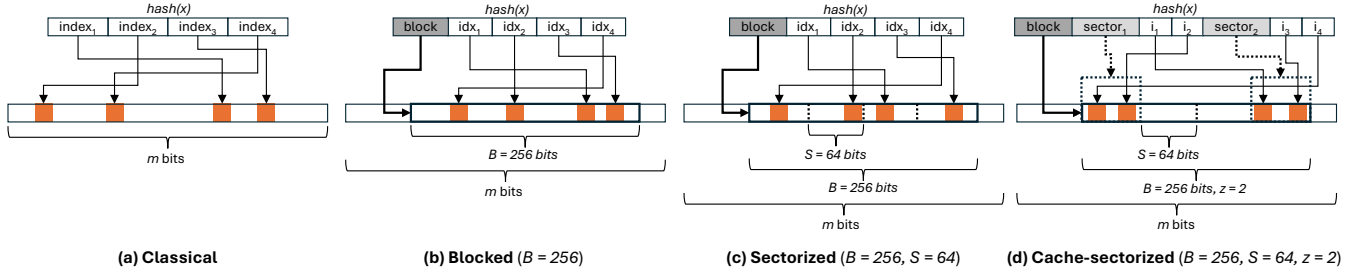Johannes Gutenberg University
Mainz, Germany

**Figure 1: Overview of Bloom filter variants. Inspired by Schmidt et al. [30]. For more details refer to Section 2.1.**

## Abstract

Bloom filters are a fundamental data structure for approximate membership queries, with applications ranging from data analytics to databases and genomics. Several variants have been proposed to accommodate parallel architectures. GPUs, with massive thread-level parallelism and high-bandwidth memory, are a natural fit for accelerating these Bloom filter variants potentially to billions of operations per second. Although CPU-optimized implementations have been well studied, GPU designs remain underexplored. We close this gap by exploring the design space on GPUs along three dimensions: vectorization, thread cooperation, and compute latency.

Our evaluation shows that the combination of these optimization points strongly affects throughput, with the largest gains achieved when the filter fits within the GPU's cache domain. We examine how the hardware responds to different parameter configurations and relate these observations to measured performance trends. Crucially, our optimized design overcomes the conventional trade-off between speed and precision, delivering the throughput typically restricted to high-error variants while maintaining the superior accuracy of high-precision configurations. At iso error rate, the proposed method outperforms the state-of-the-art by 11.35× (15.4×) for bulk filter lookup (construction), respectively, achieving above 92% of the practical speed-of-light across a wide range of configurations on a B200 GPU. We propose a modular CUDA/C++ implementation, which will be openly available soon.

## Keywords

Bloom Filter, Data Structures, Hashing, GPUs, CUDA, Vectorization, Parallel algorithms, Data Analytics

## 1 Introduction

The rapid growth of data volumes poses critical challenges for large-scale analytics. To address this, modern data processing pipelines increasingly rely on parallel architectures, such as multi-core CPUs and GPUs [2, 17], to sustain throughput and scalability. Within these pipelines, filtering plays a critical role; it eliminates irrelevant or redundant records early in the computation, thereby reducing memory pressure, communication overhead, and unnecessary downstream processing [32, 36].

A Bloom filter [4] is a space-efficient probabilistic data structure for approximate membership testing over a set of distinct elements. By design, it trades exactness for storage efficiency: membership queries may yield false positives but never false negatives. This asymmetric error property makes Bloom filters particularly desirable in applications where occasional false alarms are tolerable, but missed detections must be avoided. Prominent use cases include large-scale analysis pipelines [25], network caching [5], database systems [12, 16], and genomics applications [18, 20, 22–24, 26, 33–35].

A substantial body of research has focused on optimizing Bloom filter implementations for modern CPU architectures [1, 12, 16, 27–30]. In contrast, comparatively few efforts have explored GPU-based designs [12–14], and existing work lacks comprehensive tailoring of the implementation towards the GPU hardware architecture. This gap is striking, given that modern CUDA-enabled accelerators can provide orders of magnitude higher thread-level parallelism and memory bandwidth than CPUs. Existing GPU work either

---

*Both authors contributed equally to this research.

emphasizes extended functionality beyond the simple add/lookup semantics of Bloom filters, leading to suboptimal performance when only the basic operations are required, or lacks depth in tailoring the implementation to the specifics of the GPU architecture, leaving performance potential on the table [11, 19, 21]. Leveraging these critical architectural capabilities is key to enabling Bloom filters as a practical building block in GPU-accelerated applications.

In this work, we investigate the design and optimization of Bloom filters for modern GPUs. We consider several Bloom filter variants that are known to perform well on CPU architectures and explore how their design choices translate to GPU architectures. To tailor these designs to GPU hardware, we apply a set of optimizations that target key architectural characteristics of modern GPUs. In particular, our implementation emphasizes efficient memory access patterns through careful load and store vectorization (Section 4.1), minimizing pressure on the memory subsystem by reducing the number of issued memory operations. For fingerprint generation, we employ branchless multiplicative hashing [9] (Section 4.2), which lowers per-key latency by leveraging instruction-level parallelism while preserving sufficient fingerprint diversity. In addition, we introduce a warp-cooperative execution scheme (Section 4.3) that avoids redundant hash computations when threads jointly perform filter insertions or lookups. We evaluate our approach against a state-of-the-art GPU implementation, as well as a CPU baseline. Through a systematic exploration of the hyperparameter space, we identify configurations with favorable trade-offs between lookup and construction throughput, false-positive rate, and space consumption, and relate the observed trends to profiling results and hardware response.

Our evaluation spans both cache-resident and DRAM-resident regimes. In the latter case, the achievable throughput is fundamentally bounded by the GPU's random-access DRAM performance, providing a practical upper limit. Our implementation reaches more than 92% of this speed-of-light bound on an NVIDIA B200 and sustains this level across a broad range of filter configurations. In contrast, the competing implementation attains similar throughput only for a single configuration and at the cost of substantially higher error rates. We further show that this applies to different GPU architectures (B200, H200 SXM, and RTX PRO 6000). When the filter fits entirely within the GPU's L2 cache, the performance gap widens significantly. At comparable error rates, our approach outperforms the state-of-the-art by up to 15.4× for bulk lookups and 11.35× for filter construction, highlighting the benefits of our design in compute-dominated regimes. These contributions ensure that approximate membership testing scales seamlessly with the massive bandwidth of emerging accelerator architectures, effectively removing a critical bottleneck in large-scale data processing. Consequently, our modular library implementation can yield significant acceleration for a variety of applications, ranging from data analytics to bioinformatics.

## 2 Background

## 2.1 Bloom Filter

### 2.1.1 Classical Bloom Filter (CBF).
A Classical Bloom filter (CBF) [4] (Figure 1(a)) represents a set as a bit array of size $m$, into which $n$ distinct elements are inserted. The ratio $c = \frac{m}{n}$ denotes the number of bits allocated per element and serves as a key parameter in determining the filter's accuracy.

To insert an element, $k$ independent hash functions map it to $k$ positions in the array, and the corresponding bits are set to one. To test membership, the same $k$ positions are checked. If any bit is zero, the element is certainly not in the set; if all are one, membership is reported with a probability $f$ of being false positive caused by hash collisions.

The probability of a false positive can be approximated by

$$f = \left(1 - e^{-\frac{kn}{m}}\right)^k, \tag{1}$$

which depends on $m$, $n$, and $k$. The value of $k$ that minimizes this probability is

$$k = \frac{m}{n} \ln 2 = c \ln 2, \tag{2}$$

yielding the minimum false positive rate

$$f_{\min} = \left(\frac{1}{2}\right)^{c \ln 2}. \tag{3}$$

### 2.1.2 Blocked Bloom Filter (BBF).
While the CBF offers excellent space efficiency, its access pattern is highly inefficient on modern hardware. Each query or insertion typically reaches multiple, potentially distant positions in the bit array, leading to poor cache locality and reduced throughput. Blocked Bloom Filters (BBFs) [29] (Figure 1(b)) address this limitation by partitioning the bit array into $b$ fixed-size blocks, each of size $B = \frac{m}{b}$, where $B$ is typically chosen to match the cache line size of the underlying system. For each element, all $k$ fingerprint bits are restricted to a single block, which is selected using an additional hash function. Each block thus acts as a smaller CBF storing only $\frac{n}{b}$ items of the set, assuming the additional hash distributes elements uniformly. This organization significantly improves cache locality, since only one cache line needs to be accessed per input key, thereby reducing pressure on the memory system. The false positive rate is notably higher than that of a CBF of equal size, as the block-selection step introduces variance in the distribution of items across blocks, resulting in uneven per-block fill rates [29]. In practice, however, this accuracy penalty is overruled by the performance gains, and it can be partially offset by modestly increasing the value of $k$.

### 2.1.3 Register Blocked Bloom Filter (RBBF).
An extreme instance of the BBF is the Register Block Bloom Filter (RBBF) [16], which sets the block size equal to the machine-word size. This configuration greatly reduces computational overhead: all $k$ bits can be tested in a single word-level comparison, rather than iterating over multiple positions in a larger block, and only one word needs to be accessed per query. The drawback is a substantially higher false positive rate, since the limited number of $k$-bit patterns within a 32- or 64-bit word increases the chance of collisions.

### 2.1.4 Sectorized Bloom Filter (SBF).
Relaxing the restriction of one word per block in the RBBF leads to the Sectorized Bloom Filter (SBF) [1, 16] (Figure 1(c)), which subdivides each block into $s$ sectors of size $S = \frac{B}{s}$, where $S$ is typically chosen to match the machine-word size (e.g., 32 or 64 bits). For each key, the $k$ bits are distributed evenly across the $s$ sectors of a block, balancing bit usage and reducing intra-block contention, thereby mitigating the

error-rate limitations of the RBBF. This organization retains the efficiency of whole-word comparison operations while improving accuracy. Furthermore, because sectors are laid out contiguously in memory and can be processed independently, the SBF supports vectorized execution on both CPUs and GPUs.

*2.1.5 Cache-Sectorized Bloom Filter (CSBF).* One limitation of the SBF is the requirement that $k \geq s$, with best results when $k$ is a multiple of $s$ to ensure uniform sector contention. This constraint can force exceptionally large values of $k$ to meet error-rate requirements. For example, with $B = 1024$ and $S = 64$ (bits), we obtain $s = 16$, which already implies $k \geq 16$ – a prohibitively high value in many practical applications. Cache-Sectorized Bloom Filters (CSBF) [16] (Figure 1(d)) address this issue by adding another hierarchical layer that partitions the $s$ sectors of a block into $z$ groups. Within each group, exactly one sector is selected to store the key's $\frac{k}{z}$ fingerprint bits. Consequently, $k$ only needs to be a multiple of $z$ (with $z \leq s$), thereby allowing smaller values of $k$, and fewer words per block are updated atomically during insertion or read during lookup. The trade-off is that the sector-selection step introduces a non-uniform runtime-dependent execution path. Nevertheless, the CSBF has demonstrated superior performance in many CPU-based scenarios [16, 30]. We therefore incorporate the CSBF into our framework as a high-performance baseline to compare with the SBF, BBF, and CBF.

## 2.2 GPU memory system

Modern CUDA-enabled accelerators employ either GDDR or high-bandwidth on-chip memory (HBM). While the former typically achieve just above one terabyte per second of throughput, e.g., 1.6 TB/s on NVIDIA RTX Pro 6000 Blackwell Server Edition, HBM is capable of sustaining several terabytes per second of throughput in the optimal case, e.g., 8 TB/s (3.3 TB/s) on NVIDIA B200 (H200 SXM), respectively. This bandwidth, however, comes at the cost of relatively high access latency. The memory hierarchy consists of a two-level cache system: each Streaming Multiprocessor (SM) is equipped with a private L1 cache, while a unified L2 cache is shared across all SMs. All traffic between L1 and DRAM passes through L2. The minimum access granularity is 32B, referred to as a sector, with four sectors composing a 128B cache line. On HBM systems, or if ECC is enabled, the effective access granularity is 64B. Accesses from threads on the same SM that target the same cache line or even sector (and thus the same L2 slice) can be merged into a single L2 request through temporal coalescing. Consequently, coalesced memory access patterns are critical for efficient use of GPU DRAM, as they minimize the number of high-latency memory transactions. Importantly, atomic updates benefit from the same coalescing mechanism as well, which is essential for supporting concurrent, lock-free insertions into the filter. In the context of blocked or sectorized Bloom filters, access patterns are typically irregular: each insertion/query may target a different filter block, leading to semi-random memory accesses. Within a single filter block, accesses are contiguous and can often be coalesced into one or a few sector requests. However, when multiple blocks are processed concurrently by an SM, requests may be scattered across distant memory locations. More generally, the memory subsystem can only handle a limited number of concurrent requests. When many

outstanding accesses accumulate, either because individual requests take long to complete or because too many are issued in parallel, contention arises. In this case, new requests cannot be serviced immediately, and the issuing warp must stall until resources free up. This makes the choice of filter block size critical, as it determines the trade-off between memory access efficiency, i.e., how well each sector transfer is utilized, and per-query computational cost, defined by how many words each thread must traverse to process a single query. To avoid confusion with the CUDA terminology introduced above, we will not use the term *sector* for partitions of a SBF filter block. Instead, we refer to these partitions as *words*, reflecting the fact that the implementation maps them directly onto machine-word units which can be updated atomically.

## 3 Related Work

Early adaptations of Bloom filters for GPUs primarily focused on porting the classical variant (CBF) to the massive thread parallelism of graphics processors. Ma et al. [19] and Shi et al. [31] presented some of the first GPU implementations, utilizing global memory for the bit array. However, their approach suffered from excessive uncoalesced memory accesses inherent to the CBF design. More recent efforts have adopted blocked designs to exploit the GPU's cache hierarchy. The *WarpCore* library [14] provides a state-of-the-art BBF implementation. While efficient for specific configurations, WarpCore relies on a rigid thread-cooperation scheme that assigns a fixed number of threads to perform an insert/lookup operation. This lack of flexibility leads to suboptimal resource utilization, as the static mapping cannot adapt to the changing arithmetic intensity or register pressure of different filter configurations. Hayashikawa et al. [13] proposed the Folded Bloom Filter to improve the false-positive rate of BBFs on GPUs by distributing elements across paired bit arrays to equalize load. However, this technique necessitates complex address folding and hash generation logic that incurs computational overhead. Gubner et al. [12] explored SBF designs on GPUs within a co-processing framework for accelerating database joins. While they demonstrated the viability of the SBF layout, their work focused on the system-level orchestration of data transfers between CPU and GPU rather than the micro-architectural optimization of the filter kernel itself.

Beyond Bloom filters, several other Approximate Membership Query (AMQ) data structures have been ported to GPUs, most notably Cuckoo filters [10] and Quotient filters [3, 11]. Geil et al. [11] presented a GPU implementation of Quotient filters, which support deletions and merging. However, their design relies on linear probing and complex metadata management, resulting in significantly lower update rates compared to Bloom filters. Geil et al. reported a 2.5× slowdown for incremental insertions compared to a standard GPU Bloom filter. Similarly, GPU Cuckoo filters, as explored by McCoy et al. [21], offer higher space efficiency and deletion support. Yet, they suffer from stochastic control flow during insertion (due to cuckoo kicking sequences) and dependent memory loads during lookup (checking alternate buckets), which serializes execution and hides memory latency less effectively than the independent, parallelizable accesses of a Bloom filter. While these structures are superior when deletion support is strictly required, they cannot

match the raw throughput of an optimized SBF/CSBF for append-only or static set membership scenarios, which remain one of the dominant use cases in large-scale analytics.

## 4 Optimization Strategies

While the SBF is designed for high performance, this comes at the cost of reduced accuracy, exhibiting a much higher false-positive rate than the CBF. Peak performance is typically achieved with small block sizes ($B$), with the RBBF representing the fastest, but least accurate, extreme. Conversely, larger block sizes improve accuracy but reduce throughput. This trade-off is acceptable for some applications but prohibitive for others, motivating the use of alternative filter variants that favor accuracy over speed. In this section, we introduce several optimizations that preserve the high throughput of small-block configurations while extending performance scalability to larger block sizes. As a result, the SBF can achieve significantly improved error rates, making it competitive again in scenarios where both high performance and low false-positive rates are required.

### 4.1 Horizontal & Vertical Vectorization

Section 2.1.4 describes two key properties of SBF that are particularly relevant for efficient parallelization in terms of two dimensions:

(1) **Vertical ($\Phi$)**: Words within a block are laid out contiguously in memory, and

(2) **Horizontal ($\Theta$)**: Words can be processed independently.

When performing filter lookups, property (1) enables vectorized memory accesses, where multiple neighboring words are fetched using a single wide load instruction. On older CUDA hardware, the widest available load is 16 bytes (`ld.global.v4.u32` using the Parallel Thread Execution (PTX) ISA[1]), corresponding to half a memory sector. The Blackwell architecture introduced even wider, 32B instructions, e.g., `ld.global.v8.u32`), populating 8 32-bit registers at once. A chunk of words can thus be loaded into registers with fewer instructions and subsequently processed in a tight, statically unrolled loop while the next chunk is being prefetched concurrently, fostering instruction-level parallelism. This technique is commonly referred to as *pipelining* or *vertical vectorization* [27]. In contrast, insertions cannot directly benefit from load vectorization, since each word must be updated atomically. Nevertheless, instruction-level parallelism can still be exploited by overlapping the computation of key patterns for different words with the issuance of the corresponding atomic updates.

Property (2) enables *horizontal vectorization* [27], i.e., processing words in parallel across multiple threads. CUDA employs a Single-Instruction Multiple-Threads (SIMT) execution model, where threads are organized into warps of 32 that execute in lock-step. Warps can be further partitioned into sub-warp tiles, also known as *Cooperative Groups* (CGs), which allow finer control over the degree of parallelism. Within this setup, a CG can collectively process the words of a filter block, with each thread responsible for a subset of the work. Intra-group communication is handled by efficient

warp-collective intrinsics, e.g., register shuffles or warp-voting operations, which allow threads within a warp to exchange data directly through registers, avoiding shared or global memory access and minimizing latency. Together, vertical and horizontal vectorization provide two degrees of freedom that enable fine-grained tuning of the implementation to the characteristics of the CUDA architecture.

Let $\Phi$ denote the vertical vectorization parameter, i.e., the number of contiguous words each thread processes per step, and let $\Theta$ denote the horizontal vectorization parameter, i.e., the size of the CG that processes a filter block concurrently.

As an illustrative example, assume a block size of $B = 256$ bits and a word size of $S = 32$ bits, which yields $s = B/S = 8$ words per block, denoted $w_1, \ldots, w_s$. Figure 2 shows five representative vectorization layouts:

- **(a):** $\Theta = 1, \Phi = 8$ Corresponds to a 1:1 mapping between CUDA threads and queries. A single thread processes the words of the block sequentially. However, $\Phi = 8$ allows for loading multiple words at once, i.e., a single `ld.global.v8.u32` on Blackwell+ or two `ld.global.v4.u32` instructions back-to-back on older architectures are sufficient to fetch the entire filter block into registers, after which the words are processed in an unrolled loop.

- **(b):** $\Theta = 1, \Phi = 1$ In contrast to (a) this layout has a vertical layout of 1, meaning that a thread can only load one word at a time via `ld.global.u32`, i.e., no vectorized memory accesses are possible.

- **(c):** $\Theta = 2, \Phi = 2$ Two threads cooperatively process the block, each responsible for half of the words. With $\Phi = 2$, each thread can fetch two contiguous words for every `ld.global.v2.u32`. The block is traversed in strides of $\Theta \cdot \Phi = 4$.

- **(d):** $\Theta = 2, \Phi = 4$ Similar to (c), but each thread now covers four words at once with a single `ld.global.v4.u32`, eliminating the need for a (statically unrolled) strided loop.

- **(e):** $\Theta = 4, \Phi = 2$ Four threads cooperate on the block, each fetching two words with `ld.global.v2.u32`.

Note the constraint $1 \leq \Theta \cdot \Phi \leq s$ and both $\Theta$ and $\Phi$ must be a power-of-two.

To enforce the compiler to emit the widest possible load instructions along the $\Phi$ dimension, we employ the helper shown in Listing 1. This function loads a chunk of $\Phi$ consecutive words from the filter and returns them as a `cuda::std::array`, which the compiler can reliably place into thread-local registers. Only in cases of excessive register pressure will these values be spilled to local memory. The intrinsic `__builtin_assume_aligned` asserts that the pointer `ptr` satisfies the required alignment. At the call site, we uphold this guarantee by allocating the filter array with the proper alignment and by accessing words exclusively in strides of $\Phi$.

The parameter $\Theta$ directly determines work granularity. Larger CGs devote more threads to a single query, proportionally reducing the number of concurrent queries processed per warp. In contrast, $\Phi$ determines the per-thread workload, directly impacting the kernel's register pressure. Due to aggressive loop unrolling, along the $\Phi$ axis, a thread might require more registers than available leading to either (or both) reduced occupancy and register spilling, both of which impact performance negatively.

---

[1] NVIDIA PTX ISA Documentation https://docs.nvidia.com/cuda/parallel-thread-execution/index.html?highlight=ld#data-movement-and-conversion-instructions-ld
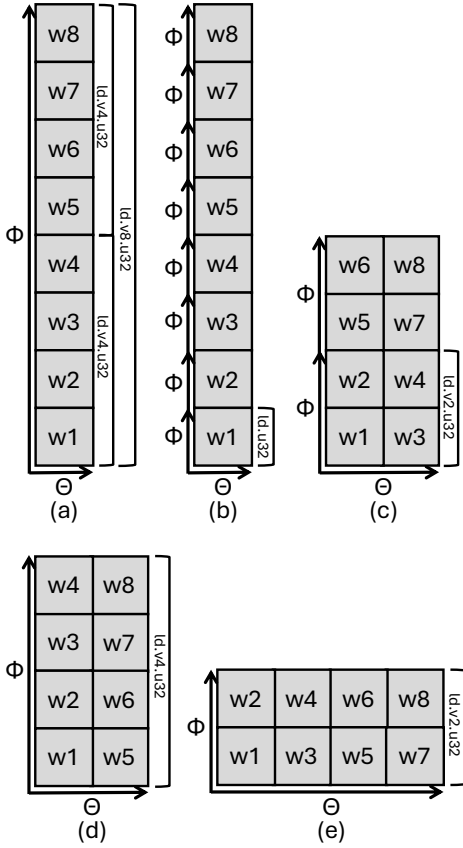
**Figure 2: Illustration of five exemplary vectorization layouts for a block size of $B = 256$ bits and a word size of $S = 32$ bits, which yields $s = B/S = 8$ words per block. The arrangement of labels $w_1, \ldots, w_8$ shows the order in which words are assigned and processed, with execution progressing from bottom to top. For each layout, the annotation on the right indicates the granularity of the memory load instruction (32, 64, 128, 256 bits) applied per step. Repeated vertical arrows denote strided processing in increments of $\Theta \cdot \Phi$. For a more detailed explanation refer to Section 4.1.**

## 4.2 Key Pattern Generation

An essential design aspect of Bloom filters is the generation of bit patterns from input keys. Typically, one or more hash functions are applied to each key to produce high-entropy hash bits that determine the bit positions to be set in the filter. A uniform distribution of these hash values minimizes collisions and thus improves accuracy. The conventional approach [4] employs $k$ independent hash functions, each producing one bit position. While this guarantees excellent pattern diversity, it incurs a computational cost that scales linearly with $k$. An alternative is to use a single hash function iteratively: the hash of the key is computed once, and subsequent hash values are derived by reapplying the same function to the key in combination with the previous hash value and an additional seed. This approach has been applied in GPU Bloom filters before [14],

```
template <size_t Phi, class Word>
__device__ cuda::std::array<Word, Phi>
    vec_load_words(Word* ptr) {
  // Maximum alignment is 32 bytes (LDG.256)
  constexpr auto alignment = min(sizeof(Word) *
      Phi, 32);
  // Provide alignment guarantee to the compiler
      to enforce vectorized load
  return *reinterpret_cast<cuda::std::array<Word,
      Phi>*>(
      __builtin_assume_aligned(ptr, alignment));
}
```

but it still requires $k$ sequential evaluations of the hash function, which limits instruction-level and thread-level parallelism. Double hashing [15] reduces the number of full hash evaluations to two and generates additional hash values through linear combinations. Although this improves efficiency, it still involves iterative computation. Another scheme, proposed by Putze et al. [29], generates fingerprints by looking up pre-computed values in a small table. While effective on some architectures, this approach proves inefficient on GPUs. Even in the optimal case where the table resides in shared memory, the random access pattern required for lookup induces significant serialization due to shared memory bank conflicts. Furthermore, the inherent latency of writing to and reading from shared memory outweighs the benefits of avoiding computation, especially when compared to arithmetic-heavy alternatives that can be hidden by the GPU's massive thread-level parallelism. Multiplicative hashing generates key patterns by multiplying a base hash value with a set of odd constants (salts). This technique has long been employed in CPU Bloom filter implementations [1, 28], and it offers two primary advantages. First, the base hash function is evaluated only once per key; all subsequent bit positions can be derived from simple integer multiplications. Second, multiplicative hashing yields universal hash functions [9], which ensure approximately uniform distribution of bits across sectors and thereby stabilize the false-positive rate. In vectorized CPU implementations, multiplicative hashing has been combined with both pure horizontal vectorization [1] and pure vertical vectorization [28]. In both cases, the number of multipliers (salts) is fixed to the number of SIMD lanes so that they fit within a single SIMD register. While this approach is sufficient for CPUs, efficient GPU execution requires greater flexibility: both horizontal and vertical vectorization (see Section 4.1) must be supported simultaneously, and the implementation must handle arbitrary combinations of $k$, $\Theta$, and $\Phi$.

Two challenges arise in this context:

(1) **Inlining salts into generated code.** A natural design is to place the multipliers in shared memory or constant memory for memory-efficient reuse during key pattern generation. However, this still requires a memory load (either `ld.shared` or `ld.const` in PTX), a synchronization barrier, and loads into registers. To eliminate the loads and reduce

access latency, the multipliers could instead be placed directly in registers, which scales poorly as $k$ increases due to increased register pressure, resulting in reduced occupancy or local memory spilling. We therefore inject the multipliers directly into the generated machine code. Using C++ template metaprogramming, we statically unroll (1) the inner loop over $\Phi$ words per thread and (2) the outer loop over $\Theta \cdot \Phi$ word groups per cooperative group. Unrolling exposes as compile-time constants both the multiplier index, which allows the compiler to inline the corresponding multiplier literal, and the outer loop index so that the compiler can narrow the range of starting multiplier indices for the threads cooperating in pattern generation in an iteration of the inner loop.

(2) **Mapping thread indices to salt indices.** When $\Theta > 1$, threads in a cooperative group must select distinct starting multipliers. A naïve solution—such as a large switch statement—would generate excessive comparisons and long predicated instruction sequences. Instead, we construct a compile-time binary decision tree over the narrowed range of starting multiplier indices. This reduces the runtime selection depth to $\log(\Theta)$, ensuring minimal branching and predication overhead.

Together, these techniques yield a mechanism that automatically specializes key pattern generation for arbitrary vectorization parameters, while maintaining inlined multipliers and efficient thread-to-salt mapping.

By default, our pattern generation method uses the 64-bit implementation of the *xxHash* algorithm [6], which offers high throughput and produces hash values with excellent entropy characteristics.

## 4.3 Adaptive Thread Cooperation

If $\Theta > 1$, the filter operation is executed in a group-cooperative manner, where two or more coalesced threads jointly process the same input key. This execution scheme requires that all participating threads operate on the same input data (key), which implies that computations independent of the group-local thread index yield the same result across the group. Such *uniform* computations are therefore redundantly executed by every thread. Although CUDA can optimize warp-wide uniform code paths through the use of special uniform registers shared across all lanes of a warp, as well as a set of warp-uniform instructions, this mechanism does not extend to sub-warp groups. As a result, with $\Theta$ cooperating threads, the number of instructions issued for otherwise uniform computations increases by a factor of $\Theta$. In the context of SBF, the computation of a key's hash value and corresponding filter block index are strictly group-uniform and thus always redundantly evaluated under this execution model. To eliminate this inefficiency, we propose an adaptive thread cooperation scheme, illustrated in Figure 3.

(1) A warp is assigned a contiguous chunk of 32 consecutive input keys, creating a 1:1 mapping between threads and inputs. This ensures coalesced memory accesses and high cache-line utilization. Each thread computes the hash value of its key and stores it in a register. For fixed-size key types and data-independent hash functions, this computation can be executed fully in lock-step (SIMD fashion) within the warp, maximizing resource utilization. (2) After this
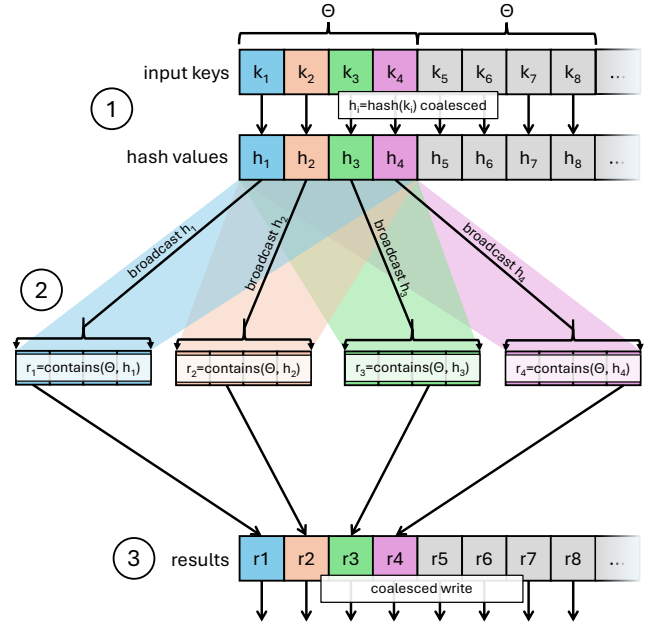


**Figure 3: Adaptive thread cooperation for a bulk filter lookup with $\Theta = 4$. The top array represents the input key sequence. Each CUDA thread is initially assigned one key and computes its hash value, storing it in a thread-local register. Subsequently, groups of $\Theta$ consecutive threads process their keys iteratively: in each iteration, the hash value of the active key is broadcast via a register shuffle, after which the group cooperatively performs the filter lookup. Each thread maintains a register for the result of its assigned key. Once all keys have been processed, the results are written back in a coalesced manner. For details, see Section 4.3.**

step, the granularity shifts from per-thread to group-cooperative execution with CG size $\Theta$. Each thread holds the hash value for one input key. The group iterates over its lanes, broadcasting the hash value of the currently active key via a register shuffle operation. All threads in the group then cooperatively execute the filter operation for that key. (3) For lookup operations, each thread stores the result for its corresponding key in a register. After all group keys are processed, the results are written back in a coalesced fashion, further improving memory efficiency.

## 5 Experimental Evaluation

We evaluate our implementation with respect to both performance and accuracy. To assess performance, we correlate the measured throughput with the practical limits of the GPU memory subsystem, providing insight into how closely the implementation approaches the hardware's empirical speed-of-light performance. We evaluate two separate scenarios due to their distinct performance characteristics: Section 5.3 focuses on filter sizes that fit into the GPU's L2 cache domain, while Section 5.2 examines filter sizes well above the cache capacity. We compare our approach against an available GPU Bloom filter implementation from the *WarpCore* [14] library (WC

BBF), a GPU CSBF integrated into our optimization framework, as well as a multithreaded CPU CSBF implementation [30]. WarpCore uses a BBF design with a static vectorization layout of $\Phi = 1$ and $\Theta$ equal to the number of machine words in the block, similar to the SBF definition. However, in contrast to an SBF design, the $k$ fingerprint bits of a key are not necessarily distributed evenly across the words, making it a BBF implementation. The GPU CSBF, to ensure a fair comparison, fully leverages our flexible vectorization scheme, multiplicative hashing, adaptive thread cooperation, and statically unrolled loops. To minimize inter-thread dependencies, we impose the following constraint on the vectorization layout: the number of words per group must be smaller than $\Phi$, the vertical vectorization parameter. The group index is calculated by introducing another odd multiplier to generate a new hash. For each valid configuration of the number of groups $z$, we perform a grid search to identify the vectorization layout satisfying this constraint that maximizes throughput.

The CPU baseline provided by Schmidt et al. [30] employs a SBF design that leverages AVX-512 vector instructions to parallelize hash computation and bit-testing. Furthermore, it applies radix partitioning to confine random memory accesses to the CPU's cache hierarchy, thereby minimizing TLB and cache misses for large filters.

## 5.1 Experimental Setup

All experiments were conducted using NVIDIA's *nvbench* [8] benchmarking library to ensure methodological consistency and were verified through profiling using NVIDIA's *Nsight Compute* [7]. The setup procedure included GPU clock locking, cache clearing between runs, thermal stabilization to prevent throttling, and repeated execution until the measurement variance fell below a predefined threshold. All experiments were performed using *CUDA Toolkit* v13.0.48 with driver v580.65 and host compiler *GCC* v13.3.0. Experiments in Sections 5.2 and 5.3 were conducted on a NVIDIA B200 GPU (Blackwell generation). CPU benchmarks were conducted on AMD EPYC 9124 16-Core. For architectural comparison (Section 5.4), additional evaluations were performed on an *NVIDIA H200 SXM* and an *NVIDIA RTX PRO 6000 Blackwell Server Edition*; the latter features GDDR7 memory instead of the on-chip HBM used in the data center–class devices.

We keep the following parameters $S = 64$, and $k = 16$ constant throughout the experiments, where the latter is required to drive the maximum tested SBF block size of $B = 1024$. For each throughput measurement run, we generate $N = 10^9$ unique, random uint64_t input keys. To benchmark the contains operation, we pre-populate the filter with these keys, ensuring that all lookups yield true positive results. Note that the key distribution does not affect throughput in this benchmark setup. The false-positive rate is measured by first inserting the space–error-rate–optimal number of distinct keys into the filter, obtained by solving Equation (3) for $n$. We then query $N$ keys not present in the insertion set and record the fraction of false-positive responses.

## 5.2 DRAM Filter

In case the filter size exceeds the GPU's L2 cache capacity, operation throughput is predominantly constrained by the random

**Table 1: Throughput (G Elem/s) of various vectorization layouts for bulk contains and add operations with a 1GB filter on B200. Best performing layouts are shown in bold. For a given value of $\Theta$ we select the maximum possible value of $\Phi$.**
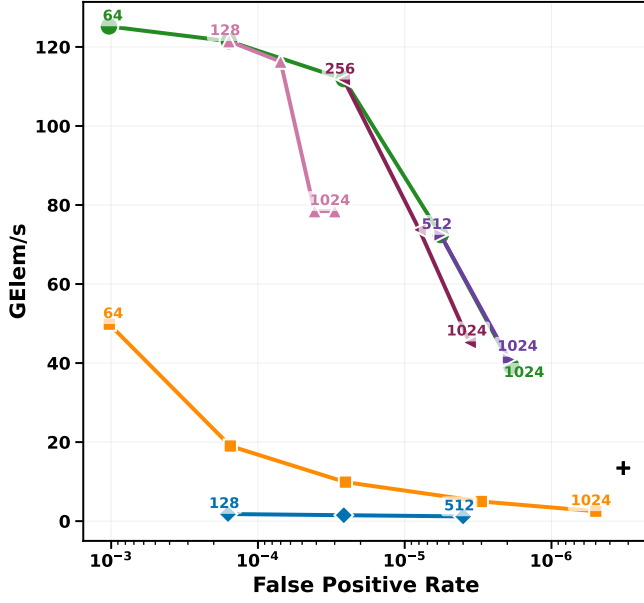
| Op. | B | $\Theta$ | | | | |
|---|---|---|---|---|---|---|
| | | 1 | 2 | 4 | 8 | 16 |
| contains | 64 | **48.69** | | | | |
| | 128 | **48.54** | 44.62 | | | |
| | 256 | **47.79** | 43.74 | 41.64 | | |
| | 512 | 25.35 | **40.66** | 40.15 | 33.66 | |
| | 1024 | 12.81 | 36.01 | **36.96** | 33.38 | 24.54 |
| add | 64 | **22.43** | | | | |
| | 128 | 13.57 | **22.26** | | | |
| | 256 | 7.59 | 13.65 | **22.10** | | |
| | 512 | 4.58 | 7.72 | 15.31 | **20.75** | |
| | 1024 | 2.88 | 5.02 | 8.53 | 15.41 | **15.61** |

access performance of the DRAM which can be quantified empirically through microbenchmarking the number of random 64-bit load/stores in giga-updates per Second (GUPS) [2]. While this setup closely resembles the access pattern of an RBBF, it also serves as an upper-bound value for larger filter block sizes.
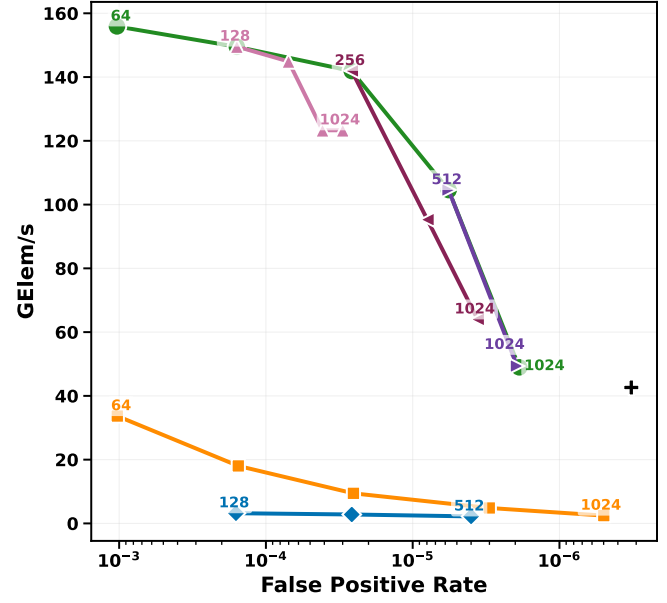
Figure 4(c) and (d) shows the throughput–error-rate frontier of several Bloom filter variants alongside the practical speed-of-light for random memory access on the B200 accelerator. Our SBF implementation achieves above 92% of the empirical speed-of-light for both construction and lookup when $B \leq 256$. Notably, reducing the block size below 256 bits does not yield additional performance gains, contrary to earlier assumptions that RBBFs ($B = 64$) offer the highest throughput. As discussed in Section 2.2, the minimum DRAM access granularity on modern GPUs is a 256-bit sector. Consequently, block sizes smaller than 256 bits do not reduce the amount of memory traffic issued to DRAM, and therefore cannot achieve higher throughput than the $B = 256$ configuration. For bulk filter lookups, the optimal vectorization layout is $\Theta = 1$ and $\Phi = s$ for $B \leq 256$. For larger block sizes, a horizontal layout of $\Theta > 1$ becomes beneficial, i.e., for $B = 512$ (1024) $\Theta = 2$ (4) is 1.6× (2.9×) faster compared to a fully vertical layout, respectively. It is beneficial to spread out the computation over multiple threads for larger filter block sizes as computational work as well as register pressure (affecting occupancy) scales proportionally to $s$. However, when using more than one thread per work item, the additional communication and synchronization among CG threads incurs additional overhead, resulting in sub-linear performance gains. Given the empirical numbers shown in Table 1, the optimal vectorization layout for lookup operations (contains) can be expressed as $\hat{\Theta}_c = max(1, \frac{B}{256})$, i.e., one thread per sector.

The performance characteristics of filter insertions (add) differ notably from lookups. Across all evaluated configurations, a fully horizontal layout consistently yielded the best performance, i.e., $\hat{\Theta}_a = s$. Profiling indicates that insertion throughput is dominated by the effectiveness of the L1 request coalescing mechanism, which
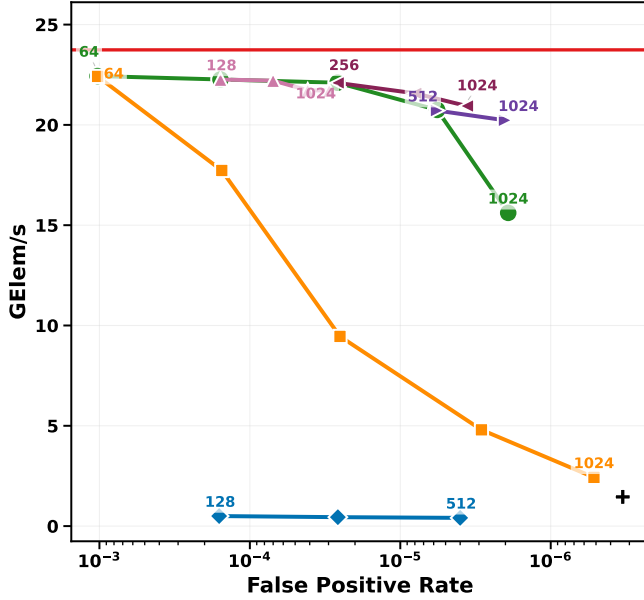
---

[2]HPC Challenge RandomAccess benchmark: http://www.hpcchallenge.org/projectsfiles/hpcc/RandomAccess.html. CUDA version: https://github.com/NVIDIA-developer-blog/code-samples/tree/master/posts/gups
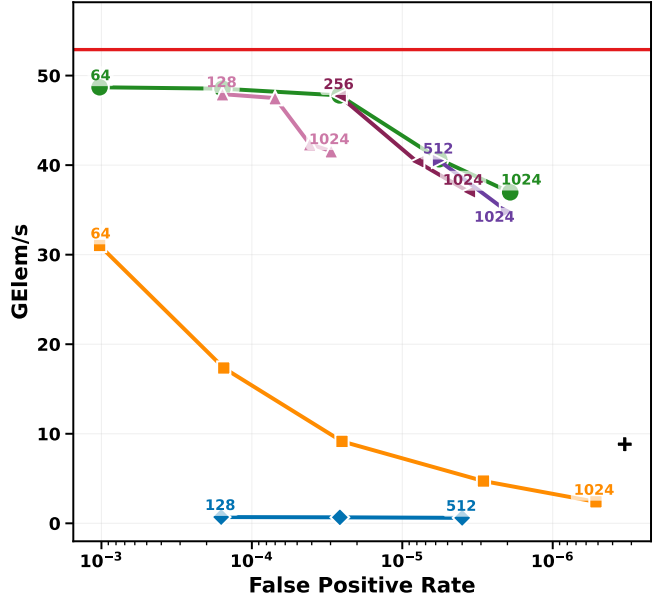
(a) Bulk construction (32 MB filter).

(b) Bulk lookup (32 MB filter).

(c) Bulk construction (1 GB filter).

(d) Bulk lookup (1 GB filter).

Figure 4: Throughput vs. false-positive rate frontier on NVIDIA B200. The top row ((a), (b)) shows performance for a 32 MB (L2-resident) filter, while the bottom row ((c), (d)) shows a 1 GB (DRAM-resident) filter. Data points are annotated with their corresponding block size $B$ in bits. The solid red line in the bottom row represents the practical speed-of-light (SOL) limit for random memory accesses.

merges per-SM accesses to the same cache line into a single request to L2. Increasing $\Theta$ from 1 to $\hat{\Theta}_a$ reduces the number of sector requests proportionally to $\frac{B}{S}$. The coalescing mechanism is temporal. Flushes occur if accesses span too many distinct cache lines over time. Thus, it is essential that all atomic updates to a given filter block are issued within the tightest possible time window. A fully horizontal layout maximizes this temporal locality: all threads targeting the same block ideally submit their atomic operations within the same cycle, enabling the coalescer to combine them into the fewest possible L2 transactions. Profiling revealed another noteworthy effect: for block sizes $B > 256$, i.e., when a lookup or insertion spans more than one sector, we observe substantial warp stalls, reported as *stall_mmio_throttle* for contains and *stall_drain* for add. These stalls arise when the memory subsystem becomes saturated with outstanding requests, preventing the issuing warp from submitting new ones. As a result, additional memory operations are delayed until sufficient capacity becomes available.

The frontiers for the GPU CSBF at larger block sizes exhibit the trade-off characteristic of this filter architecture: a smaller value of $z$, the number of groups, reduces memory traffic (benefiting memory-bound regimes) at the cost of a higher false-positive rate. However, the resulting performance gains are attenuated by the characteristics of the memory subsystem. First, the minimum DRAM access granularity of 256 bits limits bandwidth savings: when $B \geq 512$, the $z = 2$ configuration accesses two memory sectors, whereas $z \geq 4$ (and the SBF) access all four sectors of the same cache line. Second, the high latency of DRAM accesses often masks the reduction in transfer volume. This effect stands in contrast to cache-resident regimes (see Section 5.3 and Figure 4(a) and (b)), where the low latency of L2 access allows the reduced sector access count to translate more directly into throughput gains.

WarpCore's BBF implementation (WC BBF) reaches near–speed-of-light throughput for $B = 64$, but its performance declines rapidly as the block size increases. WC BBF employs a fully horizontal vectorization layout for both construction and lookup. While this layout is well suited for construction, the BBF organization induces an uneven distribution of work across words, reducing the likelihood that L1 can coalesce word updates into a single L2 transaction. For lookups, the fully horizontal layout becomes increasingly suboptimal at larger block sizes, as reflected in Table 1.

The GPU CBF baseline achieves 1.45 and 8.84 billion operations per second for insertions and lookups, respectively. Although the CBF exhibits a false-positive rate roughly two orders of magnitude lower than our $B = 256$ configuration, our implementation is 15.3× faster for add and 5.4× faster for contains. This comparison underscores the fundamental trade-off between throughput and accuracy.

The CPU SBF baseline sustains throughputs of approximately 0.45 GElem/s for construction and 0.65 GElem/s for lookups.

## 5.3 Cache-resident Filter

When the filter fits entirely within the GPU's L2 cache (Figure 4 (a) and (b)), the workload is no longer dominated by DRAM random-access throughput. Instead, profiling reveals increasing utilization of the compute pipelines as the filter block size grows. In this cache-resident regime, the performance advantages of our SBF

**Table 2: Throughput (G Elem/s) of various vectorization layouts for bulk contains and add operations with a 32MB (L2-resident) filter on B200. Best performing layouts are shown in bold. For a given value of $\Theta$ we select the maximum possible value of $\Phi$.**

| Op. | $B$ | $\Theta$ | | | | |
|---|---|---|---|---|---|---|
| | | **1** | **2** | **4** | **8** | **16** |
| contains | 64 | **155.89** | | | | |
| | 128 | **149.50** | 51.58 | | | |
| | 256 | **141.88** | 51.57 | 50.40 | | |
| | 512 | **104.55** | 50.20 | 50.35 | 45.34 | |
| | 1024 | 44.87 | **48.95** | 48.69 | 45.22 | 42.11 |
| add | 64 | **125.19** | | | | |
| | 128 | 66.07 | **121.45** | | | |
| | 256 | 33.91 | 63.25 | **111.88** | | |
| | 512 | 17.10 | 20.67 | 35.56 | **72.41** | |
| | 1024 | 8.19 | 10.37 | 11.55 | 18.91 | **39.22** |

and CSBF implementations become substantially more pronounced. Across all tested configurations, our approach exhibits markedly lower compute-pipeline congestion than WC BBF, largely due to the key-pattern generation optimizations introduced in Section 4.2. Table 2 summarizes the throughput achieved by different vectorization layouts. As in the DRAM-resident case, a fully horizontal layout provides the highest performance for filter construction. For lookups, however, a different pattern emerges: when $B \leq 512$, a purely vertical layout is substantially more effective. With $\Theta = 1$, no inter-thread communication or synchronization is required. Once $\Theta > 1$, cooperative groups must exchange data and synchronize, introducing overhead that outweighs the benefits of horizontal parallelism for these block sizes.

For the GPU CSBF, as discussed in Section 5.2, performance in the cache-resident case is more strongly affected by the sector access count (two for $z = 2$ and four for $z \geq 4$), as the lower latency of the L2 cache exposes the cost of memory traffic more visibly. Though more compute-bound in this regime, the CSBF and SBF exhibit similar instruction counts for a given block size, and the effect of the runtime-dependent group index calculation for the CSBF is invisible as the compiler can freely reorder the index calculation and the lookup on account of the static loop unrolling. With approximate parity in computation, the reduction in requests from the L2 sector for smaller $z$ becomes the deciding factor, producing a proportional improvement in throughput.

For the RBBF configuration ($B = 64$), our SBF achieves a 2.51× (4.63×) speedup over WC BBF for add (contains), while for $B = 256$, the speedup increases to 11.35× (15.4×), respectively.

The GPU CBF implementation achieves 13.43 billion insertions per second and 42.64 billion lookups per second. Interestingly, its lookup throughput exceeds that of WC BBF across all tested configurations, pointing to fundamental inefficiencies in WarpCore's kernel design for small, cache-resident filters.

Finally, a CPU SBF on a modern 16-core processor achieves throughputs in the range of 1.2 (8.8) GElem/s for construction (lookup), respectively.
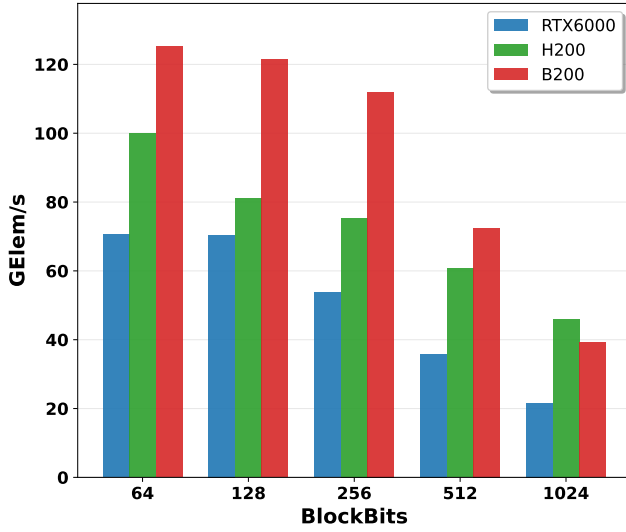
Figure 5: Bulk construction throughput of a 32MB SBF filter across various GPU architectures.
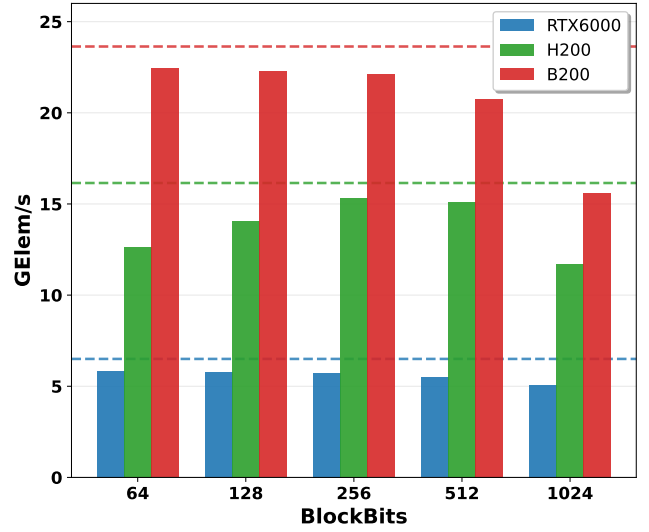


Figure 7: Bulk construction throughput of a 1GB SBF filter across various GPU architectures. Dashed lines represent the upper-bound random access throughput for a given architecture.
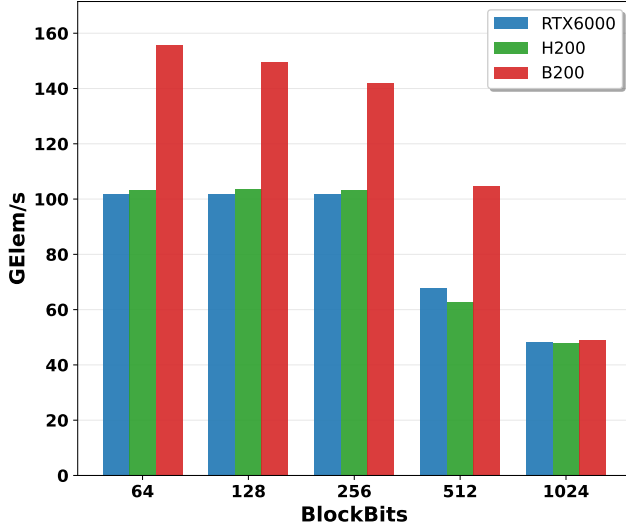


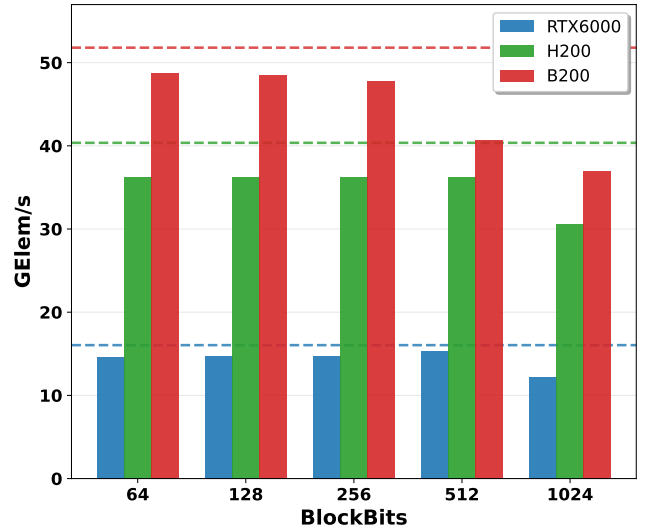Figure 6: Bulk lookup throughput of a 32MB SBF filter across various GPU architectures.



Figure 8: Bulk lookup throughput of a 1GB SBF filter across various GPU architectures. Dashed lines represent the upper-bound random access throughput for a given architecture.

## 5.4 GPU Architecture Comparison

The previous section focused on performance results obtained on NVIDIA's B200 accelerator (Blackwell generation), which features 148 SMs and ultra–high-bandwidth HBM3e memory delivering up to 8 TB/s. To better understand how our approach generalizes across GPU architectures, we extend the evaluation to two additional platforms. Specifically, we include NVIDIA's H200 SXM accelerator (Hopper generation), the predecessor of the B200, which provides 132 SMs and HBM3e memory with a substantially lower peak bandwidth of 3.3 TB/s due to fewer memory stacks. In addition, we

evaluate an NVIDIA RTX PRO 6000 Blackwell Server Edition GPU with 188 SMs. Unlike the data-center-class B200 and H200, the RTX PRO 6000 is equipped with GDDR7 memory, offering up to 1.8 TB/s of bandwidth.

The results are shown in Figures 5 to 8. For the cache-resident filter scenario, we make several interesting observations. First, the RTX PRO 6000 with its slower GDDR7 memory bandwidth is surprisingly competitive with the H200. As discussed in Section 5.3,
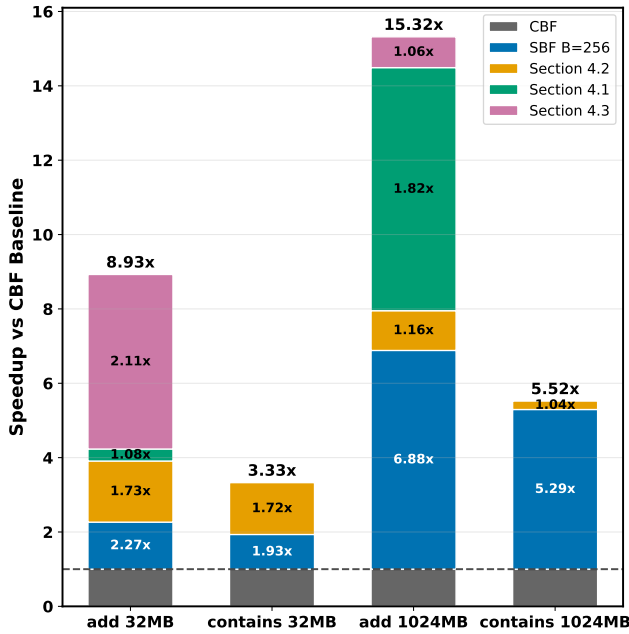
**Figure 9: Relative performance breakdown comparing our optimized SBF with a GPU CBF baseline and an unoptimized SBF ($B = 256$).**

## 5.5 Optimization Breakdown

Figure 9 summarizes the relative speedups obtained by transitioning from a GPU CBF baseline to an unoptimized SBF ($B = 256$) and incrementally applying the optimizations introduced in Section 4. Moving from a CBF to an SBF yields an immediate performance gain, which is most pronounced for DRAM-resident filters due to the reduction in random sector accesses by a factor of $k$.

Branchless multiplicative hashing (Section 4.2) primarily reduces computational overhead and has the strongest effect in the cache-resident regime, where it delivers a 1.72× speedup over the SBF baseline. Horizontal vectorization (Section 4.1) and adaptive group cooperation (Section 4.3) are only effective when $\Theta > 1$ and therefore apply exclusively to add, which performs best with a fully horizontal layout (cf. Tables 1 and 2). For contains, the optimal configuration remains $\Theta = 1$. Combined, these two optimizations yield up to a 2× improvement, with their relative impact depending on filter residency. Adaptive group cooperation is most beneficial for cache-resident filters, where low L2 latency makes the kernel increasingly compute-bound, while horizontal vectorization dominates in the DRAM-resident regime by enabling more effective coalescing of word updates into fewer sector transactions.

## 6 Conclusion

In this work, we presented a comprehensive optimization study of Bloom filters for modern GPU architectures. A central contribution of this work is the formalization of a parametric design space defined by vertical and horizontal vectorization ($\Phi, \Theta$). While applied here to membership testing, this abstraction generalizes to a wider class of irregular data structures on SIMT architectures. We showed that decoupling the logical vectorization layout from the physical thread organization allows the implementation to adapt to the distinct bottlenecks of the cache-resident and DRAM-resident regimes. We integrated this flexible vectorization with branchless multiplicative hashing and an adaptive thread cooperation scheme. This combination mitigates redundant work during group-cooperative execution and minimizes instruction overhead.

In the DRAM-resident regime, our implementation saturates the memory subsystem, achieving over 92% of the theoretical random-access bandwidth for filters with block sizes $B \leq 256$. This confirms that aligning filter block sizes with the GPU's minimum sector access granularity is critical for maximizing throughput. In the cache-resident regime, where execution becomes compute-bound, our optimizations yield substantial gains. We achieve up to 15.4× higher lookup throughput and 11.35× higher construction throughput compared to the state-of-the-art WarpCore implementation at comparable error rates. Crucially, this performance allows our implementation to deliver the throughput typically reserved for high-error Register Blocked Bloom Filters (RBBF), but with the superior accuracy of large-block variants. The level of efficiency achieved—sustained across diverse configurations—unlocks new throughput capabilities for next-generation data analytics, genomics, and database systems running on high-bandwidth accelerators.

To facilitate reproducibility and adoption, we will publicly release our modular CUDA/C++ implementation under open-source license.

L2-resident workloads are not purely bound by random-access DRAM throughput but become increasingly compute-bound with larger filter block sizes due to the low latency of the L2 cache. The RTX PRO 6000's 42% SM advantage over the H200 (188 vs. 132 SMs) and its newer Blackwell architecture allow it to exploit this compute-bound regime effectively. For the add operation on L2-resident filters, all three architectures achieve similar peak throughput at their respective optimal configurations, though H200 exhibits a preference for lower horizontal vectorization ($\Theta = 4$ at $B = 512$, $\Theta = 8$ at $B = 1024$) compared to B200 and RTX PRO 6000 (which both prefer $\Theta = 8$ and $\Theta = 16$, respectively), possibly reflecting differences in L2 cache organization or atomic operation handling between Hopper and Blackwell. For the DRAM-resident filter scenario, throughput differences are more pronounced and correlate strongly with each platform's random-access memory bandwidth. Using GUPS microbenchmarks to establish speed-of-light bounds, we measure 52.9/23.7 GUPS (read/write) for B200, 40.4/16.2 GUPS for H200, and 16.0/6.5 GUPS for RTX PRO 6000. Across all three architectures, our implementation achieves ~90–95% of these bounds: B200 reaches 92% (read) and 95% (write), H200 achieves 90% and 95%, while RTX PRO 6000 attains 95% and 90%, respectively. This high efficiency validates that our vectorization approach effectively saturates random-access memory bandwidth regardless of memory technology (HBM3e vs. GDDR7). Another notable architectural divergence appears in the contains operation: RTX PRO 6000 prefers higher horizontal vectorization ($\Theta = 2$) at intermediate block sizes ($B = 128, 256$) where both data-center GPUs favor $\Theta = 1$.

**Future Work.** We plan to explore asynchronous memory operations, specifically the global-to-shared asynchronous copy instructions (`cp.async` and `cp.async.bulk`; the latter utilizing the Tensor Memory Accelerator (TMA) unit available on Hopper+), to better overlap memory latency with hash computations and block comparison operations. However, our preliminary experiments suggest that the overhead incurred by the roundtrip through shared memory currently negates potential performance gains. We intend to conduct a deeper analysis to determine if this bottleneck can be overcome.

# References

[1] Jim Apple. 2021. Split block Bloom filters. *CoRR* abs/2101.01719 (2021). arXiv:2101.01719 https://arxiv.org/abs/2101.01719

[2] Felipe Aramburú, William Malpica, Kaouther Abrougui, Amin Aramoon, Romulo Auccapuclla, Claude Brisson, Matthijs Brobbel, Colby Farrell, Pradeep Garigipati, Joost Hoozemans, Supun Kamburugamuve, Akhil Nair, Alexander Ocsa, Johan Peltenburg, Rubén Quesada López, Deepak Sihag, Ahmet Uyar, Dhruv Vats, Michael Wendt, Jignesh M. Patel, and Rodrigo Aramburú. 2025. Theseus: A Distributed and Scalable GPU-Accelerated Query Processing Platform Optimized for Efficient Data Movement. *CoRR* abs/2508.05029 (2025). arXiv:2508.05029 doi:10.48550/ARXIV.2508.05029

[3] Michael A. Bender, Martin Farach-Colton, Rob Johnson, Russell Kraner, Bradley C. Kuszmaul, Dzejla Medjedovic, Pablo Montes, Pradeep Shetty, Richard P. Spillane, and Erez Zadok. 2012. Don't Thrash: How to Cache Your Hash on Flash. *Proc. VLDB Endow.* 5, 11 (2012), 1627–1637. doi:10.14778/2350229.2350275

[4] Burton H. Bloom. 1970. Space/Time Trade-offs in Hash Coding with Allowable Errors. *Commun. ACM* 13, 7 (1970), 422–426. doi:10.1145/362686.362692

[5] Andrei Z. Broder and Michael Mitzenmacher. 2003. Survey: Network Applications of Bloom Filters: A Survey. *Internet Math.* 1, 4 (2003), 485–509. doi:10.1080/15427951.2004.10129096

[6] Yann Collet. 2012. *xxHash: Extremely fast non-cryptographic hash algorithm.* Cyan4973. https://github.com/Cyan4973/xxHash Accessed: 2025-10-04.

[7] NVIDIA Corporation. 2025. *Nsight Compute - NVIDIA CUDA Kernel Profiling Tool.* https://developer.nvidia.com/nsight-compute Accessed: 2025-12-09.

[8] NVIDIA Corporation. 2025. *nvbench: CUDA Kernel Benchmarking Library.* https://github.com/NVIDIA/nvbench Accessed: 2025-10-05.

[9] Martin Dietzfelbinger, Torben Hagerup, Jyrki Katajainen, and Martti Penttonen. 1997. A Reliable Randomized Algorithm for the Closest-Pair Problem. *J. Algorithms* 25, 1 (1997), 19–51. doi:10.1006/JAGM.1997.0873

[10] Bin Fan, David G. Andersen, Michael Kaminsky, and Michael Mitzenmacher. 2014. Cuckoo Filter: Practically Better Than Bloom. In *Proceedings of the 10th ACM International on Conference on emerging Networking Experiments and Technologies, CoNEXT 2014, Sydney, Australia, December 2-5, 2014,* Aruna Seneviratne, Christophe Diot, Jim Kurose, Augustin Chaintreau, and Luigi Rizzo (Eds.). ACM, 75–88. doi:10.1145/2674005.2674994

[11] Afton Geil, Martin Farach-Colton, and John D. Owens. 2018. Quotient Filters: Approximate Membership Queries on the GPU. In *2018 IEEE International Parallel and Distributed Processing Symposium, IPDPS 2018, Vancouver, BC, Canada, May 21-25, 2018.* IEEE Computer Society, 451–462. doi:10.1109/IPDPS.2018.00055

[12] Tim Gubner, Diego G. Tomé, Harald Lang, and Peter A. Boncz. 2019. Fluid Co-processing: GPU Bloom-filters for CPU Joins. In *Proceedings of the 15th International Workshop on Data Management on New Hardware, DaMoN 2019, Amsterdam, The Netherlands, 1 July 2019,* Thomas Neumann and Ken Salem (Eds.). ACM, 9:1–9:10. doi:10.1145/3329785.3329934

[13] Masatoshi Hayashikawa, Koji Nakano, Yasuaki Ito, and Ryota Yasudo. 2019. Folded Bloom Filter for High Bandwidth Memory, with GPU Implementations. In *2019 Seventh International Symposium on Computing and Networking, CANDAR 2019, Nagasaki, Japan, November 25-28, 2019.* IEEE, 18–27. doi:10.1109/CANDAR.2019.00011

[14] Daniel Jünger, Robin Kobus, André Müller, Christian Hundt, Kai Xu, Weiguo Liu, and Bertil Schmidt. 2020. WarpCore: A Library for fast Hash Tables on GPUs. In *27th IEEE International Conference on High Performance Computing, Data, and Analytics, HiPC 2020, Pune, India, December 16-19, 2020.* IEEE, 11–20. doi:10.1109/HIPC50609.2020.00015

[15] Adam Kirsch and Michael Mitzenmacher. 2008. Less hashing, same performance: Building a better Bloom filter. *Random Struct. Algorithms* 33, 2 (2008), 187–218. doi:10.1002/RSA.20208

[16] Harald Lang, Thomas Neumann, Alfons Kemper, and Peter A. Boncz. 2019. Performance-Optimal Filtering: Bloom overtakes Cuckoo at High-Throughput. *Proc. VLDB Endow.* 12, 5 (2019), 502–515. doi:10.14778/3303753.3303757

[17] Viktor Leis, Peter Boncz, Alfons Kemper, and Thomas Neumann. 2014. Morsel-driven parallelism: a NUMA-aware query evaluation framework for the many-core age. In *International Conference on Management of Data, SIGMOD 2014, Snowbird, UT, USA, June 22-27, 2014,* Curtis E. Dyreson, Feifei Li, and M. Tamer Özsu (Eds.). ACM, 743–754. doi:10.1145/2588555.2610507

[18] Yiming Li, Haoling Zhang, Yuxin Chen, Yue Shen, and Zhi Ping. 2024. DNA Bloom filter enables anti-contamination and file version control for DNA-based data storage. *Briefings Bioinform.* 25, 3 (2024). doi:10.1093/BIB/BBAE125

[19] Lin Ma, Roger D. Chamberlain, Jeremy D. Buhler, and Mark A. Franklin. 2011. Bloom Filter Performance on Graphics Engines. In *International Conference on Parallel Processing, ICPP 2011, Taipei, Taiwan, September 13-16, 2011,* Guang R. Gao and Yu-Chee Tseng (Eds.). IEEE Computer Society, 522–531. doi:10.1109/ICPP.2011.27

[20] Ketil Malde and Bryan O'Sullivan. 2009. Using Bloom Filters for Large Scale Gene Sequence Analysis in Haskell. In *Practical Aspects of Declarative Languages, 11th International Symposium, PADL 2009, Savannah, GA, USA, January 19-20, 2009. Proceedings (Lecture Notes in Computer Science, Vol. 5418),* Andy Gill and Terrance Swift (Eds.). Springer, 183–194. doi:10.1007/978-3-540-92995-6_13

[21] Hunter McCoy, Steven A. Hofmeyr, Katherine A. Yelick, and Prashant Pandey. 2023. High-Performance Filters for GPUs. In *Proceedings of the 28th ACM SIGPLAN Annual Symposium on Principles and Practice of Parallel Programming, PPoPP 2023, Montreal, QC, Canada, 25 February 2023 - 1 March 2023,* Maryam Mehri Dehnavi, Milind Kulkarni, and Sriram Krishnamoorthy (Eds.). ACM, 160–173. doi:10.1145/3572848.3577507

[22] Nathaniel McVicar, Chih-Ching Lin, and Scott Hauck. 2017. K-Mer Counting Using Bloom Filters with an FPGA-Attached HMC. In *25th IEEE Annual International Symposium on Field-Programmable Custom Computing Machines, FCCM 2017, Napa, CA, USA, April 30 - May 2, 2017.* IEEE Computer Society, 203–210. doi:10.1109/FCCM.2017.23

[23] Svenja Mehringer, Enrico Seiler, Felix Droop, Mitra Darvish, René Rahn, Martin Vingron, and Knut Reinert. 2023. Hierarchical Interleaved Bloom Filter: enabling ultrafast, approximate sequence queries. *Genome Biology* 24, 1 (2023), 131. doi:10.1186/s13059-023-02971-4

[24] Páll Melsted and Jonathan K. Pritchard. 2011. Efficient counting of k-mers in DNA sequences using a Bloom Filter. *BMC Bioinform.* 12 (2011), 333. doi:10.1186/1471-2105-12-333

[25] Ripon Patgiri, Sabuzima Nayak, and Samir Kumar Borgohain. 2019. Role of Bloom Filter in Big Data Research: A Survey. *CoRR* abs/1903.06565 (2019). arXiv:1903.06565 http://arxiv.org/abs/1903.06565

[26] David Pellow, Darya Filippova, and Carl Kingsford. 2017. Improving Bloom Filter Performance on Sequence Data Using $k$-mer Bloom Filters. *J. Comput. Biol.* 24, 6 (2017), 547–557. doi:10.1089/CMB.2016.0155

[27] Orestis Polychroniou, Arun Raghavan, and Kenneth A. Ross. 2015. Rethinking SIMD Vectorization for In-Memory Databases. In *Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data, Melbourne, Victoria, Australia, May 31 - June 4, 2015,* Timos K. Sellis, Susan B. Davidson, and Zachary G. Ives (Eds.). ACM, 1493–1508. doi:10.1145/2723372.2747645

[28] Orestis Polychroniou and Kenneth A. Ross. 2014. Vectorized Bloom filters for advanced SIMD processors. In *Proceedings of the Tenth International Workshop on Data Management on New Hardware* (Snowbird, Utah, USA) *(DaMoN '14).* Association for Computing Machinery, New York, NY, USA, Article 6, 6 pages. doi:10.1145/2619228.2619234

[29] Felix Putze, Peter Sanders, and Johannes Singler. 2009. Cache-, hash-, and space-efficient bloom filters. *ACM J. Exp. Algorithmics* 14 (2009). doi:10.1145/1498698.1594230

[30] Tobias Schmidt, Maximilian Bandle, and Jana Giceva. 2021. A four-dimensional analysis of partitioned approximate filters. *Proceedings of the VLDB Endowment* 14, 11 (2021), 2355–2368. doi:10.14778/3476249.3476286

[31] Haixiang Shi, Bertil Schmidt, Weiguo Liu, and Wolfgang Müller-Wittig. 2010. A parallel algorithm for error correction in high-throughput short-read data on CUDA-enabled graphics hardware. *Journal of Computational Biology* 17, 4 (2010), 603–615. doi:10.1089/CMB.2009.0062

[32] Lakshmikant Shrinivas, Sreenath Bodagala, Ramakrishna Varadarajan, Ariel Cary, Vivek Bharathan, and Chuck Bear. 2013. Materialization strategies in the Vertica analytic database: Lessons learned. In *29th IEEE International Conference on Data Engineering, ICDE 2013, Brisbane, Australia, April 8-12, 2013,* Christian S. Jensen, Christopher M. Jermaine, and Xiaofang Zhou (Eds.). IEEE Computer Society, 1196–1207. doi:10.1109/ICDE.2013.6544909

[33] Sanjay Kumar Srikakulam, Sebastian Keller, Fawaz Dabbaghie, Robert Bals, and Olga V. Kalinina. 2023. MetaProFi: an ultrafast chunked Bloom filter for storing and querying protein and nucleotide sequence data for accurate identification of functionally relevant genetic variants. *Bioinform.* 39, 3 (2023). doi:10.1093/BIOINFORMATICS/BTAD101

[34] Henrik Stranneheim, Max Käller, Tobias Allander, Björn Andersson, Lars Arvestad, and Joakim Lundeberg. 2010. Classification of DNA sequences using Bloom filters. *Bioinform.* 26, 13 (2010), 1595–1600. doi:10.1093/BIOINFORMATICS/BTQ230

[35] Deyou Tang, Yucheng Li, Daqiang Tan, Juan Fu, Yelei Tang, Jiabin Lin, Rong Zhao, Hongli Du, and Zhongming Zhao. 2022. KCOSS: an ultra-fast k-mer counter for assembled genome analysis. *Bioinform.* 38, 4 (2022), 933–940. doi:10.1093/BIOINFORMATICS/BTAB797

[36] Yifei Yang, Hangdong Zhao, Xiangyao Yu, and Paraschos Koutris. 2024. Predicate Transfer: Efficient Pre-Filtering on Multi-Join Queries. In *14th Conference on Innovative Data Systems Research, CIDR 2024, Chaminade, HI, USA, January 14-17, 2024.* www.cidrdb.org. https://www.cidrdb.org/cidr2024/papers/p22-yang.pdf